

# Determining Association Between Categorical Variables

9

WEDNESDAY

2022 FEBRUARY

2022 JANUARY

S	T	M	T	W	T	F	S
30	31						1
2	3	4	5	6	7	8	
9	10	11	12	13	14	15	
16	17	18	19	20	21	22	
23	24	25	26	27	28	29	

WEEK 07  
040-325

8

Bad Avg. Very good Excellent

→ Ordinal data

Variable = "Food Quality"

## Determining and Inferring association

### Conditional Probability

- Ex:- Consider B-school which shortlisted 1200 candidates (960 men and 240 women) for its post graduate management program. Out of there, 324 candidates were given offer letters for admission.
- The data is included here:

	Male	Female	Total	
Offers made	288	36	324	Contingency
Not offered	672	204	876	Table
Total	960	240	1200	

- After reviewing the record, a women's forum raised the issue of gender discrimination on the basis that 288 male candidates were offered admission against only 36 female candidates.

Event of randomly selected male candidate =  $\frac{288}{960} = \frac{3}{10}$   
 Event of randomly selected female candidate =  $\frac{36}{960} = \frac{1}{20}$

Event of admission offer being made = A

" randomly selected candidate for admission offer being made = N

JANUARY

MARCH 2022

M	T	W	T	F	S	S
1	2	3	4	5	6	
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

THURSDAY  
FEBRUARY 2022

10

WEEK 07  
041-324

- To this, the B-school management replied that it is not a case of discrimination, but was because of the fact that only 240 women candidates appeared for the examination.

Let  $A^C$  be the event that the candidate is not offered admission.

One can see that  $\Pr(A) + \Pr(A^C) = 1$

- Probability that a randomly observed candidate is a male and is offered the admission. (Joint event)

$$\Pr(M \cap A) = 288/1200 = 0.24$$

- Probability that a randomly observed candidate is a male and is not offered the admission.

$$\Pr(M \cap A^C) = 672/1200 = 0.56$$

- Similarly,

$$\Pr(W \cap A) = 36/1200 = 0.03$$

$$\Pr(F \cap A^C) = 204/1200 = 0.17$$

Joint Probability

In terms of probabilities, the previous table can now be written as:

	Male	Female	Total
Offers made	0.24	0.03	0.27
Not offered	0.56	0.17	0.73
Total	0.8	0.2	1.0

$\square$  = Joint probability

$\square$  = Marginal probability

Marginal probability

MARCH

APRIL

Inspiron 15  
3000 Series

11

FRIDAY

2022 FEBRUARY

2022 JANUARY

S	M	T	W	T	F	S
30	31					1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

WEEK 07  
042-323

- What will be  $\Pr(A|M)$ ? (Conditional Probability)

The probability of event A happening, given the fact that event M has already taken place.

- This conditional probability tell us that we are concerned with admission status of only males!
- We know that out of the 960 male candidates, 288 were offered admission. So probability that a male candidate is offered admission will be  $288/960 = 0.3$

- Also observe that:

$$\Pr(A|M) = \frac{288}{960} = \frac{288/1200}{960/1200} = \frac{0.24}{0.8} = 0.3 = \frac{\Pr(A \cap M)}{\Pr(M)}$$

Joint probability  
Marginal Probability

$$\boxed{\text{Conditional Probability} = \frac{\text{Joint Probability}}{\text{Marginal Probability}}}$$

$$\Pr(A|W) = \frac{\Pr(A \cap W)}{\Pr(W)} = \frac{0.03}{0.2} = 0.15$$

MARCH 2022

M	T	W	T	F	S	S
1	2	3	4	5	6	
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

SATURDAY  
FEBRUARY 2022

12

WEEK 07  
043-322Conclusion

- The probability of admission offer given the candidate is a male is 0.3, twice of 0.15 probability of admission offer given by the candidate is a woman.
- Although the use of conditional probability does not, in itself, prove discrimination, there is support for the argument!

so

Baye's Rule

- Given the new information, we can update our prior beliefs by calculating revised probabilities - this is called the posterior probability.
- Baye's Rule is used to calculate the posterior probability if we have the initial belief (probability) and the additional sample information.

Ex:-

- Suppose that a manufacturer receives some raw material from two different suppliers  $S_1$  and  $S_2$ .
- Currently 65% of the raw material comes from  $S_1$  and remaining, 35%, comes from  $S_2$ .
- Also, suppose that from the historical data available with the quality assurance department, we know that  $S_1$  has 98% of the supplied raw material of

MARCH

APRIL

13

SUNDAY

14

MONDAY

2022 FEBRUARY

WEEK-08  
045-320Inspiron 15  
3000 Series

2022 JANUARY

S	M	T	W	T	F	S
30	31					1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

good quality and  $S_2$  has 95% of the raw material of good quality.

- That is, the probability of a "Good" quality raw material given that the supplier is  $S_1$  is  $Pr(G|S_1) = 0.98$ . And for the second supplier, this probability is:  $Pr(G|S_2) = 0.95$ .

Q. What is the probability of the raw material being supplied by  $S_1$  and it being good?  
 Soln:- Joint probability, of course!

According to the Baye's formula

$$\text{Conditional Probability} = \frac{\text{Joint Probability}}{\text{Marginal Probability}}$$

$$\text{Joint Prob.} = \text{Cond. Prob.} \times \text{Marg. Prob.}$$

$$Pr(G|S_1) = \frac{Pr(G \cap S_1)}{Pr(S_1)}$$

$$Pr(G \cap S_1) = Pr(G|S_1) \times Pr(S_1)$$

$$= 0.98 \times 0.65 = 0.637$$

$$Pr(G \cap S_2) = Pr(G|S_2) \times Pr(S_2)$$

$$= 0.85 \times 0.95 = 0.3325$$

JANUARY

MARCH 2022

M	T	W	T	F	S	S
1	2	3	4	5	6	
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

TUESDAY

FEBRUARY 2022

15

WEEK 08

- Now, knowing all this information so far, suppose the manufacturer inspects the incoming raw material on receipt and finds a bad quality material.

- He wants to know the supplier who needs to be contacted to complain!

10

- We are interested in the posterior probability that a particular supplier is guilty of supply bad quality product given that we have bad quality raw material at our doorstep -  $Pr(S_1|B)$  or  $Pr(S_2|B)$

- This is an application of Baye's theorem - finding posterior probability given some initial facts and numbers.

- From Baye's formula we know that :

$$\begin{aligned} Pr(S_1|B) &= \frac{Pr(S_1 \cap B)}{Pr(B)} \\ &= \frac{Pr(S_1) \times Pr(B|S_1)}{Pr(B)} \end{aligned}$$

- What is  $Pr(B)$ ?

- That is the probability of receiving a bad quality raw material.

- Now bad quality raw material can from supplies of  $S_1$  or  $S_2$ .

MARCH

APRIL

16

WEDNESDAY

2022 FEBRUARY

WEEK 08  
047-318

S	M	T	W	T	F	S
30	31					1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

- That is, the event  $B$  can occur with  $S_1$  or with  $S_2$ .

$$\Pr(B) = \Pr(S_1 \cap B) + \Pr(S_2 \cap B)$$

$$\text{But } \Pr(S_i \cap B) = \Pr(S_i) \times \Pr(B|S_i), \text{ and}$$

$$\Pr(S_2 \cap B) = \Pr(S_2) \times \Pr(B|S_2)$$

$$\begin{aligned} \Pr(S_1|B) &= \frac{\Pr(S_1) \times \Pr(B|S_1)}{\Pr(S_1) \times \Pr(B|S_1) + \Pr(S_2) \times \Pr(B|S_2)} \\ &= \frac{0.65 \times 0.02}{0.65 \times 0.02 + 0.35 \times 0.05} = 0.426 \end{aligned}$$

$$\Pr(S_2|B) = 0.574$$

- Significance: Find posterior probability using prior information
- Notice that we use  $\Pr(B|S_1)$  to find  $\Pr(S_1|B)$

# Inferential Association Between Categorical Variables: Chi-Squared Test of Independence

MARCH 2022

M	T	W	T	F	S	S
1	2	3	4	5	6	
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

THURSDAY  
FEBRUARY 2022

17

WEEK 08  
048-317

## Example: Brand preferences

- Suppose a survey is conducted in Mumbai and Chennai asking respondents their preferences about three brands. The result is summarized below:

City	Brand A	Brand B	Brand C	Total
Mumbai	279	73	225	577
Chennai	165	47	191	403
Total	444	120	416	980

- Independent (explanatory) variable is the city.
- Dependent (response) variable is the brand preference
- We know how to summarize the data by calculating the marginal and joint probabilities.
- Two categorical variables are statistically independent if the population condition distribution on one of them is identical to each category of the other.
- In the example, the two conditional distributions are not identical. e.g. Brand A is preferred more in Mumbai than in Chennai.
- Refer to the same example extended to a third city:

City	Brand A	Brand B	Brand C	Total
Mumbai	440(44%)	140(14%)	420(42%)	1000(100%)
Chennai	44(44%)	14(14%)	42(42%)	100(100%)
Delhi	110(44%)	35(14%)	105(42%)	250(100%)

18

FRIDAY

2022 FEBRUARY

WEEK 08  
049-316

2022 JANUARY

S	M	T	W	T	F	S
30	31					1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

- Conditional distributions is same across the cities. Hence we can conclude that brand preference is independent of the cities.
- However, statistical independence is a ~~summary~~ symmetric property b/w two categorical variable.
- If the conditional distributions within the rows are identical, then so are the distributions within the columns.
- One can verify that the conditional distributions amongst columns equals (74%, 13%, 19%)
- Based on this single sample information, we can draw inferences about the population, as we have been doing?

### Chi-square distribution

- Null hypothesis -  $H_0$ : The categorical variables are independent.
- Alternate hypothesis -  $H_1$ : The categorical variables are not independent.

Let  $O$  be the observed frequencies (from the sample)  
 Let  $E$  be the expected frequencies, if the variables were independent.

$$\begin{array}{cccccc} (N_{11})_{001} & (N_{12})_{01} & (N_{13})_{01} & (N_{14})_{01} & \dots & \\ (N_{21})_{025} & (N_{22})_{025} & (N_{23})_{025} & (N_{24})_{025} & \dots & \end{array}$$

JANUARY

MARCH 2022

M	T	W	T	F	S	S
1	2	3	4	5	6	
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

SATURDAY  
FEBRUARY 2022

19

WEEK 08  
050-315

The expected frequency for a cell equals the product of row and column totals for that cell, divided by total sample size.

- Calculate the expected frequency:

$$\text{Brand A (Mumbai)} = \frac{444 \times 577}{980} = 261.4$$

$$\text{Brand A (Chennai)} = \frac{444 \times 403}{980} = 182.6$$

$$\text{Brand B (Mumbai)} = \frac{120 \times 577}{980} = 70.7$$

$$\text{Brand B (Chennai)} = \frac{120 \times 403}{980} = 49.3$$

City	Brand A	Preferred brand			Total
		Brand B	Brand C		
Mumbai	279 (261.4)	73 (70.7)	225 (244.9)		577
Chennai	165 (182.6)	47 (49.3)	191 (171.1)		403
Total	444	120	416		980

Observed frequency      Expected frequency

- Chi-Squared test statistic:

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e} = \frac{(279 - 261.4)^2}{261.4} + \frac{(47 - 49.3)^2}{49.3} + \frac{(191 - 171.1)^2}{171.1} = 7.0$$

21

MONDAY

2022 FEBRUARY

WEEK 09  
052-313Inspiron 15  
3000 Series

15

2022 JANUARY						
S	M	T	W	T	F	S
30	31					1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

- When the  $H_0$  is true, expected and observed frequencies tend to be close for each cell, and the test statistic value is relatively small.
- If  $H_0$  is false, at least some cells have a big gap b/w expected and observed frequencies, leading to a large test statistic value.
- The larger the  $\chi^2$  value, greater is the evidence against the null hypothesis of independence.
- Degrees of freedom for the Chi-Squared distribution is given by the expression :  $df = (rows - 1) \times (cols - 1)$ .  

$$df = (r-1) \times (c-1)$$
. r and c are the # of rows and columns respectively.

$\chi^2 = 7.0$ ,  $df = 2$ , so at  $\alpha = 0.05$  (95% confidence), the tabular value of test statistic,  $\chi^2 = 5.99$

- So we reject the null hypothesis of independence.
- However, at  $\alpha = 0.01$  (99% confidence), the tabular value of test statistic,  $\chi^2 = 9.21$ , and we can not reject the null hypothesis.

MARCH 2022

M	T	W	T	F	S	S
1	2	3	4	5	6	
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31			

TUESDAY  
FEBRUARY 2022

22

WERK 09  
053-312

Box Plot

Outlier detection

Line Chart

Trend Analysis

Column Stacked Chart

Proportion Comparison parts of a

Histogram

To guess distribution

Scatter Plot

Correlation Analysis

Bar Chart

To do comparison

Pie Chart

Proportion

H<sub>0</sub> (null hypothesis) :- The population follows the proposed distribution.H<sub>a</sub> (alternative hypothesis) :- The population does not follow the proposed distribution

PQ. In the goodness-of-fit test, if the computed test statistic is greater than the tabulated value of the statistic at a given significance level then,

At the specified significance level, reject the null hypothesis and conclude that the data does not come from population with proposed distribution.

Q. A distribution is Left Tailed if...

- (a) Negative Skewness
- (b) Left Tail implied mean is on left side
- (c) Mean < Median < Mode

MARCH

APRIL

23

WEDNESDAY

2022 FEBRUARY

WEEK 09  
054-311

2022 JANUARY

S	M	T	W	T	F	S
30	31					1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

- Poisson distribution - Mean & Variance are same
- Exponential - Mean & Median are not same
- Standard normal dist. - Mean is 0
- Uniform distribution - Mean, Median and Mode are very close to each other

$$\rightarrow = \sigma^2$$

$$\text{Expected frequency} = \frac{\text{Total observed value}}{\text{No. of bins}}$$